

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2007-524144

(P2007-524144A)

(43) 公表日 平成19年8月23日(2007.8.23)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 9/50 (2006.01)	G06F 9/46 465A	5B042
G06F 11/30 (2006.01)	G06F 11/30 G	

審査請求 有 予備審査請求 未請求 (全 24 頁)

(21) 出願番号	特願2006-515645 (P2006-515645)	(71) 出願人	505467959
(86) (22) 出願日	平成16年3月16日 (2004.3.16)		フジツブ シーメンス コンピュータース
(85) 翻訳文提出日	平成17年12月19日 (2005.12.19)		ゲゼルシャフト ミット ベシュレンク
(86) 国際出願番号	PCT/DE2004/000530		テル ハフツング
(87) 国際公開番号	W02004/114143		ドイツ連邦共和国 80807 ミュンヘン
(87) 国際公開日	平成16年12月29日 (2004.12.29)		ドマークシュトラッセ 28
(31) 優先権主張番号	10327601.7	(74) 代理人	100075166
(32) 優先日	平成15年6月18日 (2003.6.18)		弁理士 山口 巖
(33) 優先権主張国	ドイツ (DE)	(72) 発明者	フリース、ベルンハルト
(31) 優先権主張番号	10330322.7		ドイツ連邦共和国 68766 ホッケン
(32) 優先日	平成15年7月4日 (2003.7.4)		ハイム イン デア クラム 28
(33) 優先権主張国	ドイツ (DE)	(72) 発明者	ナヴァビ、グラーナ
			ドイツ連邦共和国 76669 パート
			シェンボルン ツォイターナー シュトラッセ 35

最終頁に続く

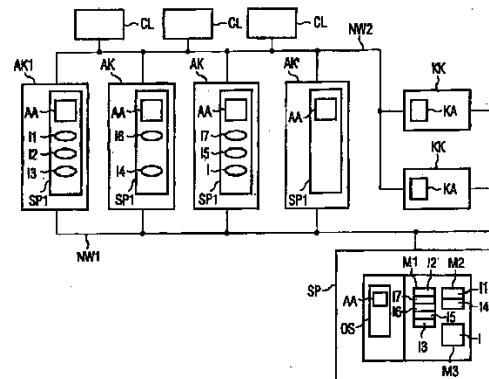
(54) 【発明の名称】 クラスタ装置

(57) 【要約】

【課題】 少ない構成費用にて作動させることのできるクラスタ装置を提供する。

【解決手段】 第1のネットワーク (NW1) と、アプリケーションノード (AK) を成しかつオペレーティングシステム (OS) を有する少なくとも2つのデータ処理システムと、各アプリケーションノード (AK) におけるアプリケーションエージェント (AA) とを有するクラスタ装置において、アプリケーションエージェント (AA) が次の機能を有する。すなわち、

- アプリケーションノード (AK) において実行されるインスタンス (I1, I2, I3) の監視 (UB) および認識
- 新たなインスタンス (I3) の始動またはアプリケーションノードにおいて予定よりも早く終了させられたインスタンスの再始動、
- アプリケーションノード (AK) における新たなインスタンス (I3) の実行が可能であるか否かの評価 (BE) および判定、
- ネットワーク (NW1) に接続されているアプ



【特許請求の範囲】

【請求項 1】

- 第 1 のネットワーク (NW1) を備えていること、
 - それぞれ 1 つのアプリケーションノード (AK) を成しかつそれぞれ 1 つのオペレーティングシステム (OS) を有する少なくとも 2 つのデータ処理システムを備え、アプリケーションノード (AK) は実行されるインスタンス (I1, I2, I3, I4) を有すること、
 - 各アプリケーションノード (AK) に、次の機能を有するアプリケーションエージェント (AA) を備えていること、すなわち、
 - アプリケーションノード (AK) において実行されるインスタンス (I1, I2, I3) の監視 (UB) および識別；
 - 新たなインスタンス (I3) の始動 (ST) またはアプリケーションノードにおいて予定よりも早く終了させられたインスタンスの再始動；
 - アプリケーションノード (AK) における新たなインスタンス (I3) の実行が可能であるか否かの評価 (BE) および決定；
 - ネットワーク (NW1) に接続されているアプリケーションノードのアプリケーションエージェント (AA) へのインスタンス実行要求 (AF)；
 - ネットワーク (NW1) に接続されているアプリケーションノード (AK) のアプリケーションエージェント (AA) へのインスタンス (I3) の実行要求 (AF) の引き受け後のメッセージ通知 (ME)；
- を特徴とするクラスタ装置。

10

20

【請求項 2】

アプリケーションエージェントの監視 (UB) の機能はリスト (T) の作成 (L) を含み、リスト (T) は、それぞれアプリケーションノード (AK) において実行されるインスタンス (I1, I2) と、実行されるインスタンスの実行に必要な全てのデータ (D) とを部分リスト (TI1, TI2) として含んでいることを特徴とする請求項 1 記載のクラスタ装置。

【請求項 3】

アプリケーションエージェント (AA) の監視 (UB) の機能は、アプリケーションノードにおいて実行されるインスタンスと他のインスタンスおよび / またはパラメータとの関連性を認識するように構成されていることを特徴とする請求項 1 又は 2 記載のクラスタ装置。

30

【請求項 4】

アプリケーションエージェント (AA) は、インスタンスの不安定な動作状態を認識するように構成されている機能を有することを特徴とする請求項 1 乃至 3 の 1 つに記載のクラスタ装置。

【請求項 5】

アプリケーションエージェント (AA) は、アプリケーションノードにおいて実行されるインスタンスを終了させるための機能を有することを特徴とする請求項 1 乃至 4 の 1 つに記載のクラスタ装置。

40

【請求項 6】

リスト (T) はアプリケーションノードに関する情報も有することを特徴とする請求項 2 記載のクラスタ装置。

【請求項 7】

記憶装置 (SP) は、第 1 のネットワークに接続されていて、かつアプリケーションノード (AK) において実行可能な少なくとも 1 つのインスタンス (I1) を有することを特徴とする請求項 1 乃至 6 の 1 つに記載のクラスタ装置。

【請求項 8】

アプリケーションノードのアプリケーションエージェントは、アプリケーションノードにおいて動作するオペレーティングシステム (OS) のサービスであることを特徴とする

50

請求項 1 乃至 7 の 1 つに記載のクラスタ装置。

【請求項 9】

記憶装置 (S P) にはアプリケーションノード (A K) のためのオペレーティングシステム (O S) が格納されていることを特徴とする請求項 7 記載のクラスタ装置。

【請求項 10】

新たなインスタンス (I 3) の実行が可能であるアプリケーションノード (A K ') が設けられていることを特徴とする請求項 1 乃至 9 の 1 つに記載のクラスタ装置。

【請求項 11】

クラスタ装置は第 1 のネットワークに接続された制御ノード (K K) として構成された少なくとも 1 つのデータ処理システムを有し、制御ノード (K K) がオペレーティングシステムおよび制御エージェント (K A) を有し、制御エージェント (K A) は、次の機能、すなわち、

- 第 1 のネットワーク (N W 1) に接続されているアプリケーションノード (A K) の機能性の検査 (U P) ;
- そのネットワーク (N W 1) に接続されているアプリケーションノード (A K) のアプリケーションエージェント (A A) へのインスタンス実行要求 (A F) ;
- アプリケーションノード (A K) の決定 (B S) およびこのアプリケーションノードへの新たなインスタンスの実行要求 ;

を有することを特徴とする請求項 1 乃至 10 の 1 つに記載のクラスタ装置。

【請求項 12】

アプリケーションノード (A K) の検査 (U P) の際に、検査すべきアプリケーションノード (A K) は、アプリケーションノード (A K) のアプリケーションエージェント (A A) によって作成されたリスト (T) を介して求められ得ることを特徴とする請求項 11 記載のクラスタ装置。

【請求項 13】

実行要求 (A F) は、アプリケーションエージェント (A A) の作成された部分リスト (T I 1 , T I 2) を有することを特徴とする請求項 2 又は 11 記載のクラスタ装置。

【請求項 14】

作成されたリスト (T I 1) および / または実行要求 (A F) および / またはメッセージ (M E) は、少なくとも 1 つのファイルとして記憶装置 (S P) に格納されていることを特徴とする請求項 1 乃至 13 の 1 つに記載のクラスタ装置。

【請求項 15】

記憶装置 (S P) に格納されている各インスタンスは、インスタンスが実行されるアプリケーションノード (A K) に割り当てられる一義的な識別番号を有することを特徴とする請求項 1 乃至 14 の 1 つに記載のクラスタ装置。

【請求項 16】

第 1 のネットワーク (N W 1) は、 T C P / I P または N F S プロトコルによる通信を行うように構成されていることを特徴とする請求項 1 乃至 15 の 1 つに記載のクラスタ装置。

【請求項 17】

少なくとも 1 つの制御ノード (K K) はアプリケーションノード (A K) の初期化プロセスを制御するための手段を有し、初期化プロセスはアプリケーションノード (A K) のオペレーティングシステム (O S) の読み込みのためのコマンドを有することを特徴とする請求項 11 乃至 16 の 1 つに記載のクラスタ装置。

【請求項 18】

少なくとも 2 つのアプリケーションノード (A K , A K ') が第 2 のネットワーク (N W 2) に接続され、第 2 のネットワーク (N W 2) はクライアントコンピュータ (C L) に接続されていることを特徴とする請求項 1 乃至 17 の 1 つに記載のクラスタ装置。

【請求項 19】

記憶装置 (S P) に格納されているインスタンスはデータベースを有することを特徴と

10

20

30

40

50

する請求項 1 乃至 18 の 1 つに記載のクラスタ装置。

【請求項 20】

少なくとも 1 つのデータ処理システムに識別番号 (IP1, IP2) が割り当てられていることを特徴とする請求項 1 乃至 19 の 1 つに記載のクラスタ装置。

【請求項 21】

クラスタ装置のデータ処理システム (DV) は同じ識別番号 (IP1, IP2) によりプール (VC1, VC2, VC3) を成すことを特徴とする請求項 20 記載のクラスタ装置。

【請求項 22】

アプリケーションノード (AK) のアプリケーションエージェント (AA) は、次の機能

- アプリケーションノード (AK) に割り当てられた識別番号 (IP1, IP2) の評価；
- 同じ識別番号 (IP1, IP2) を有するアプリケーションノード (AK) におけるアプリケーションエージェント (AA) へのインスタンスの実行要求；
- 同じ識別番号 (IP1, IP2) を有するアプリケーションノード (AK) におけるアプリケーションエージェント (AA) へのインスタンス実行要求の引き受け後のメッセージ通知；

を有することを特徴とする請求項 20 乃至 21 の 1 つに記載のクラスタ装置。

【請求項 23】

識別番号 (IP1, IP2) は IP アドレスまたは IP アドレスの一部を含むことを特徴とする請求項 20 乃至 22 の 1 つに記載のクラスタ装置。

【請求項 24】

制御ノード (KK1) として構成されたデータ処理システム (DV) に識別番号 (IP2) が割り当てられていて、制御ノード (KK1) において実行される制御エージェント (KA) は同じ識別番号を有するアプリケーションノード (AK) の機能を検査するように構成されていることを特徴とする請求項 20 乃至 23 の 1 つに記載のクラスタ装置。

【請求項 25】

アプリケーションノード (AK) において実行されるインスタンス (I1, I2, L1, L2, L3) に優先順位 (PS1, PS2) が割り当てられていて、アプリケーション

【請求項 26】

制御ノード (KK1) の制御エージェント (KA) は、評価、判定およびこれらのインスタンスの実行要求の際にアプリケーションノード (AK) において実行されるインスタンスの優先順位 (PS1, PS2) を評価するように構成されていることを特徴とする請求項 20 乃至 25 の 1 つに記載のクラスタ装置。

【請求項 27】

ネットワーク (NW1) に接続されている少なくとも 2 つのアプリケーションノード (AK, AK') および制御ノード (KK) からなるクラスタ装置における方法であって、

- 制御ノード (KK) が、アプリケーションノード (AK, AK') から、ノードにおいて実行される全てのインスタンス (I1, I2, I3) とインスタンスの実行のために必要なデータとを有するリスト (T) を受け取り、
- 制御ノード (KK) がアプリケーションノード (AK) の故障を規則的な時間間隔にて検査し、
- 制御ノード (KK) が、アプリケーションノード (AK, AK') の故障時に、故障アプリケーションノードにおいて実行されるインスタンスと実行に必要なデータとを有するリスト (TI3) を作成して、ネットワーク (NW1) に接続されているアプリケーションノード (AK, AK') に実行要求と共に転送することを特徴とする方法。

10

20

30

40

50

【請求項 28】

制御ノード（KK）が、アプリケーションノードの故障時に、アプリケーションノードによって作成されたリスト（TI3）を実行要求と一緒に少なくとも1つの他のアプリケーションノード（AK）に転送することを特徴とする請求項 27 記載の方法。

【請求項 29】

アプリケーションノードがインスタンスの実行のためのリスト（T）を作成し、このリスト（T）と共に要求を少なくとも1つの他のアプリケーションノードに伝えることを特徴とする請求項 27 記載の方法。

【請求項 30】

制御ノード（KK）が、アプリケーションノードの故障時に、アプリケーションノードによって作成されたリスト（T）と調整可能なパラメータとにより他のアプリケーションノードを確定し、これに故障アプリケーションノードにおいて実行されるインスタンスの実行要求を伝達することを特徴とする請求項 27 又は 28 記載の方法。

10

【請求項 31】

制御ノード（KK）が、検査すべきアプリケーションノードを、アプリケーションエージェントが受け取ったリスト（T）の評価によって確定することを特徴とする請求項 27 乃至 30 の1つに記載の方法。

【請求項 32】

方法が、制御ノードにおいて実行される制御エージェント（KA）およびアプリケーションノードにおいて実行されるアプリケーションエージェント（AA）によって実施されることを特徴とする請求項 27 乃至 31 の1つに記載の方法。

20

【請求項 33】

少なくとも1つのアプリケーションノード（AK）および制御ノード（KK）に識別番号（IP2）が割り当てられていて、

- 制御ノード（KK1）の制御エージェント（KA）が、同じ識別番号（IP2）を有する少なくとも1つのアプリケーションノード（AK）のアプリケーションエージェント（AA）からテストメントを受け取ることを特徴とする請求項 27 乃至 32 の1つに記載の方法。

【請求項 34】

- 実行されるインスタンス（I1, I2, L1, L2, L3）に優先順位（PS1, PS2）が割り当てられ、

30

- 制御エージェント（KA）が、アプリケーションノード（AK）の故障時に、故障アプリケーションノード（AK）において実行されるインスタンスに割り当てられた優先順位（PS1, PS2）を評価し、優先順位（PS1, PS2）に依存して実行要求を送信することを特徴とする請求項 17 乃至 30 の1つに記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ネットワークに接続されている少なくとも2つのアプリケーションノードおよび1つの制御ノードからなるクラスタ装置およびクラスタ装置における方法に関する。

40

【0002】

ネットワークを介して結合されかつ設定された課題を共通に処理する複数のコンピュータからなる相互接続はクラスタと呼ばれる。処理すべき課題は小さな部分課題に分解され、個々のコンピュータに割り振られる。この種の公知のクラスタが、とりわけ非常に計算費用のかかる課題に使用される Biowulf クラスタ (www.biowulf.org) である。他の形式のクラスタの場合には、クラスタの計算速度ではなくて可用性が関心の的である。この形式のクラスタの場合には、クラスタ内の1つのコンピュータの故障時に他のコンピュータが故障コンピュータの課題を時間損失なしにまたは僅かな時間損失のみで引き継ぐことが保証されなければならない。このようなクラスタの例がインターネット内のウェブサーバであり、またはリレーショナルデータベースを有する中央記憶装置

50

でもある。

【0003】

このような動作様式を有するクラスタは、高可用性クラスタとも呼ばれ、ネットワークを介して互いに接続されている複数の個別サーバを有する。各サーバはクラスタのノードを形成する。アプリケーションが処理されるサーバはアプリケーションノードと呼ばれ、中央の管理、制御または監視の課題を有するサーバは制御ノードを形成する。アプリケーションノードでは種々のアプリケーションまたは大きなアプリケーションの種々の部分アプリケーションが実行され、個々のアプリケーションは互いに関連することができる。クライアントと呼ばれるクラスタ外のコンピュータは、クラスタ内で実行されるアプリケーションにアクセスしてデータを呼び出す。

10

【0004】

このようなクラスタは、アプリケーションノードのほかに、中央インスタンスである制御ノードを含む。制御ノードは、個々のアプリケーションノードにおける実行中のアプリケーションを監視し、場合によってはこれを終了させ、またはそれを新たに始動させる。1つのアプリケーションノードの故障時には、中央インスタンスが残りのアプリケーションノード上において、中止されたアプリケーションを新たに始動させる。このために制御ノードは、なおも十分に容量を有するノードを選択する。その際に、クラスタの構成および稼働率に応じて、今まで使用されていなかったアプリケーションノードが使用され、または新たに始動されるアプリケーションの計算負荷ができるだけ一様に配分される。この過程が負荷バランシングと呼ばれる。

20

【0005】

他方では中央インスタンスまたは制御ノードを故障から守るために、大抵は中央インスタンスの役目をする他のサーバによって中央インスタンスに冗長性を持たせることが必要である。しかしながら、このようなクラスタ解決策は、アプリケーションノードと中央インスタンスとの間におけるデータ交換が非常に大きくなるという欠点を有する。これに加えて、各アプリケーションノードは中央インスタンスの照会に回答するために計算時間を消費する。更に、制御ノードはあらゆる有り得る故障シナリオを処理することができなければならないために、構成費用およびそれにつながる欠陥のある構成のリスクが少なからず増大する。

【0006】

本発明の課題は、明白に少ない構成費用にて作動させることのできるクラスタ装置を提供することにある。

30

【0007】

本発明は並列関係にある特許請求の範囲によって解決される。

【0008】

本発明による装置においては、第1のネットワークと、それぞれ1つのアプリケーションノードを成しかつそれぞれ1つのオペレーティングシステムを有する少なくとも2つのデータ処理システムとを備えたクラスタが設けられている。各アプリケーションノードはアプリケーションエージェントを有し、少なくとも1つのアプリケーションノードは実行されるインスタスを有する。アプリケーションエージェントは少なくとも次の機能を含んでいる。すなわち、

40

- アプリケーションノードにおいて実行される全てのインスタスの機能性および誤りのない動作の監視、
- 新たなインスタスの自立した始動またはアプリケーションノードにおいて予定よりも早く終了したインスタスの再始動、
- アプリケーションノードにおける新たなインスタスの実行が可能であるか否かの評価および決定、
- ネットワークに接続されているアプリケーションノードのアプリケーションエージェントへのインスタス実行要求、
- ネットワークに接続されているアプリケーションノードのアプリケーションエージェ

50

ントへのインスタンス実行要求の引き受け後のメッセージ通知。

【0009】

各アプリケーションノードには、自立的にかつ他のアプリケーションノードにおけるアプリケーションエージェントならびに中央インスタンスに依存せずに動作するアプリケーションエージェントまたはプログラムが設けられている。特に、アプリケーションエージェントまたはプログラムは、新たなインスタンスの始動もしくは実行が可能であるか否かに関する自立した決定到達ができるように構成されている。アプリケーションエージェントは、アプリケーションエージェントがインスタンス実行要求を受け取るか、またはインスタンスの実行が誤りのない動作のために必要になるときに、この機能を実行する。評価は記憶、要求に含まれる情報および予め定義された規則を含めて行なうことが望ましい。

10

【0010】

付加的に、アプリケーションエージェントは他のアプリケーションエージェントへの要求のための機能を有する。それによって、エージェントには、差し迫った故障の際に、自立して中央インスタンスに依存せずに、アプリケーションノードにおいて実行されるインスタンスを他のアプリケーションノードに引き渡すことが可能である。代替として、この機能により、インスタンス実行要求を他のアプリケーションエージェントに送信することもできる。この機能は、実行すべきインスタンスの誤りのない動作のための全てのパラメータを準備するという趣旨で実施されている。

20

【0011】

特に、アプリケーションエージェントはこれらの機能により、アプリケーションノードに限定された自立した負荷バランシングを行なうことができる。

【0012】

各アプリケーションノードにおけるアプリケーションエージェントの待機により中央インスタンスは負担を軽減され、中央インスタンスとアプリケーションノードとの間におけるデータ通信が低減され、監視機能がアプリケーションノードへ移される。アプリケーションエージェントはそれらの側において独立している。

【0013】

アプリケーションエージェントがインスタンス実行要求を拒絶するとき、他のアプリケーションノードのアプリケーションエージェントへのメッセージを発生する機能を有することが望ましい。

30

【0014】

アプリケーションノードにおけるアプリケーションエージェントの監視機能がリスト作成を含む場合に特に有利である。リストは、それぞれアプリケーションノードにおいて実行されるインスタンスと、実行されるインスタンスの実行に必要な全てのデータおよびパラメータを含んでいる。望ましくは、リストが部分リストに分割されていて、各部分リストが1つのインスタンスのためのデータを含んでいる。更に、これらのリストがアプリケーションノードに関する情報および動作パラメータも含んでいると有利である。これらは可能な実行に関する評価の際に使用されることが好ましい。特に簡単なやり方では、ネットワークに接続されているアプリケーションノードのアプリケーションエージェントへのリストもしくは部分リストの発送による要求が行なわれる。

40

【0015】

本発明による発展形態では、クラスタ装置は第1のネットワークに接続されている記憶装置を有する。記憶装置はアプリケーションノードのアクセスのために構成されている。記憶装置は、アプリケーションノードにおいて実行可能な少なくとも1つのインスタンスを含んでいる。それによって、クラスタ装置のアプリケーションノードは、記憶装置内のインスタンスにアクセスして実行のために取り込むことができる。全てのデータを記憶装置に保存し、これらを全てのアプリケーションノードに利用可能にすることが望ましい。この共通利用される記憶装置はコストを低減し、クラスタ装置のメンテナンスを簡単化す

50

る。

【0016】

これに関連して、記憶装置にアプリケーションノードのためのオペレーティングシステムが格納されているとよい。それによって、オペレーティングシステムは各アプリケーションノードに個々にインストールされなくてもよく、アプリケーションノードの初期化過程において記憶装置から読み込まれる。したがって、オペレーティングシステムにおける更新および変更を簡単に実施することができる。各アプリケーションノードのアプリケーションエージェントがアプリケーションノードで作動するオペレーティングシステムのサービスであることが望ましい。アプリケーションエージェントは共通利用される記憶装置に格納されている。アプリケーションエージェントがアプリケーションノードの初期化の際に自動的に始動されると有意義である。 10

【0017】

クラスタ装置の実施形態においては、新たなインスタンスの実行が可能であるアプリケーションノードが設けられている。したがって、クラスタ装置は、アプリケーションノードの故障時に、故障ノードにおいて動作中のアプリケーションを引き継ぐアプリケーションノードを常に含んでいる。

【0018】

有利な発展形態においては、クラスタ装置は制御ノードとして構成された少なくとも1つのデータ処理システムを有し、データ処理システムは第1のネットワークに接続されている。制御ノードはオペレーティングシステムおよび制御エージェントを有し、制御エージェントは、次の機能を有する。すなわち、 20

- 第1のネットワークに接続されているアプリケーションノードの機能性の検査、
- そのネットワークに接続されているアプリケーションノードのアプリケーションエージェントへのインスタンス実行要求、
- アプリケーションノードの決定およびこのアプリケーションノードへの新たなインスタンスの実行要求。

【0019】

制御ノードにおけるこのような制御エージェントを用いることにより、アプリケーションノードと制御ノードとの間のデータ交換が明白に低減される。特に、アプリケーションノードの機能性の検査を簡単な周期的な存在テストによって行なうことができる。合理的な構成においては、アプリケーションノードのアプリケーションエージェントがその存在および機能性を制御エージェントによって検査される。個々のインスタンスの検査は省略される。なぜならば、これはアプリケーションエージェントによって行なわれるからである。アプリケーションノードの全体的な故障の際には制御エージェントが故障インスタンスの実行のための新たなアプリケーションノードを決定する。それによって常に誤りのない動作が保証される。 30

【0020】

アプリケーションノードの検査の際に、検査すべきアプリケーションノードを、アプリケーションノードのアプリケーションエージェントによって作成されたリストを介して求めることができる。その際にアプリケーションノードのアプリケーションエージェントによって作成されたリストを制御エージェントが任意に使用することができ、制御エージェントがこれを評価する。リストによって制御エージェントはクラスタ装置内に存在するアプリケーションノードに関する知識を得る。したがって、クラスタ装置内に存在するノードの動的な探索が不要となる。特に、本発明による装置のこの構成においては、クラスタ装置に更なるアプリケーションノードを簡単に追加することができる。制御エージェントへのリストの伝達後に、新たなアプリケーションノードの周期的な検査が行なわれる。 40

【0021】

インスタンスの実行のための要求は、アプリケーションノードのアプリケーションエージェントの作成された部分リストを有することが好ましい。実行要求の機能は、制御エージェントにおいて、またアプリケーションエージェントにおいて、等しく構成されている 50

ことが望ましい。

【0022】

本発明の望ましい発展形態においては、アプリケーションエージェントによって作成されたリストおよび/または実行要求および/または実行要求引き受けメッセージが記憶装置の少なくとも1つのファイルとして格納されている。これは、中央インスタンスもしくは制御インスタンスの存在なしに、各アプリケーションエージェントによるアクセスおよび独立した評価を可能にする。更に、要求が記憶装置のメモリ領域におけるリストの簡単な準備によって通知され、このリストの除去によって引き受け後の通報がなされることが好ましい。

【0023】

本発明の発展形態では、記憶装置に格納された各インスタンスが一義的な識別番号を有し、識別番号はインスタンスが実行されるアプリケーションノードに割り当てられる。それによって、各アプリケーションノードにおける各インスタンスを既に実行された他のインスタンスに依存せずに実行させることができる。第2のインスタンスへの第1のインスタンスのアクセスは一義的な識別番号を介して行なわれる。それによって、個々のインスタンスおよびアプリケーションエージェントはクラスタ装置の構造的な構成に関する知識を必要としない。1つのノードにおいて複数のインスタンスが実行される場合には、複数の識別番号を割り当てることも勿論可能である。一実施形態においては一義的な識別番号は仮想のIPアドレスである。

【0024】

第1のネットワークはTCP/IPまたはNFSプロトコルによる通信を行うように構成されているとよい。両プロトコルは多数のデータ伝送および管理の可能性をもたらし、特に簡単に実現することができる。

【0025】

これに関連して、少なくとも1つの制御ノードはアプリケーションノードの初期化プロセスを制御するための手段を有し、初期化プロセスはアプリケーションノードのオペレーティングシステムの読み込みのためのコマンドを有することが望ましい。したがって、この手段により、制御ノードによってアプリケーションノードにおける初期化プロセスが始動され、初期化プロセスがオペレーティングシステムのロードをもたらす。特に、初期化プロセスの間にアプリケーションノードのためのコマンドもしくはパラメータが引き渡される。

【0026】

本発明の発展形態では、クラスタ装置における少なくとも2つのアプリケーションノードおよび少なくとも1つの制御ノードが第2のネットワークに接続され、第2のネットワークがクライアントコンピュータに接続されている。このコンピュータは、大抵はアプリケーションノードにおいて実行されるインスタンスへの照会を送信するために使用される。第2のネットワークの構成によって、アプリケーションノード間のデータ流と、クライアントコンピュータとアプリケーションノードとの間のデータ流とが分離される。したがって、1つのネットワークにおける交換データ量が低減され、同時に分離によって監視または不当なアクセスに対する信頼性が高められる。

【0027】

特別に有利な構成は次の構成である。すなわち、記憶装置に格納されているインスタンスがデータベースの一部として構成されていることである。代替的には、格納されたインスタンスはデータベースにアクセスするアプリケーションである。この場合にデータベースは記憶装置の一部であることが好ましい。

【0028】

本発明の他の発展形態においては、各データ処理システムに識別番号が割り当てられている。同じ識別番号を有するクラスタ装置のデータ処理システムはプールを成す。したがって、簡単なやり方でクラスタ装置が更に分割されて、個別の課題を引き受けることができる。

10

20

30

40

50

【0029】

ネットワークに接続されている少なくとも2つのアプリケーションノードおよび制御ノードからなるクラスタ装置における方法は、制御ノードが、アプリケーションノードからアプリケーションノードにおいて実行される全てのインスタンスとインスタンスの実行のために必要なデータおよびパラメータとを有するリストを受け取ることを特徴とする。更に、制御ノードがアプリケーションノードを故障について規則的な時間間隔にて検査し、アプリケーションノードの故障時に、故障アプリケーションノードにおいて実行されるインスタンスと実行に必要なデータとを有するリストを作成する。制御ノードは、このリストを、ネットワークに接続されているアプリケーションノードに実行要求と共に転送する。

【0030】

この方法により、制御ノードはアプリケーションノードの故障を、換言するならば、アプリケーションノードの存在を検査するだけである。アプリケーションノードにおいて実行されるインスタンスの検査、管理または監視は放棄される。それによってアプリケーションノードと制御ノードとの間のデータ量が明白に低減される。本方法の発展形態では、アプリケーションノードが、ノードで実行されるインスタンスの状態変化時に、変化通知または新しいリストを送信する。

【0031】

制御ノードは、監視機能によって記録されたアプリケーションノードの故障の際に、アプリケーションノードによって作成されたリストを実行要求と一緒に少なくとも1つの他のアプリケーションノードに転送する。代替として、アプリケーションノードが、実施のために決定されたインスタンスと実施のために必要なデータとを有するリストを作成し、このリストを少なくとも1つの他のアプリケーションノードに伝える。

【0032】

他の構成においては、制御ノードが、アプリケーションノードの故障時に、アプリケーションノードによって作成されたリストと調整可能なパラメータとにより他のアプリケーションノードを求める。求められたアプリケーションノードには、故障アプリケーションノードにおいて実行されるインスタンスの実行要求が送信される。それによって、効率的なやり方にて負荷バランシングがクラスタ全体のための制御ノードにより行なわれる。これに関連して、制御ノードが、第1のアプリケーションノードに対しては、インスタンス終了のための信号を送信し、第2のアプリケーションノードに対しては、その終了させられたインスタンスを実行させるための信号を送信する。

【0033】

特に、作成されたリストおよび調整可能なパラメータにより、実行を予定よりも早く終了させられたインスタンスの実行のために適切なコンピュータを見つけ出すことができる。更に、制御ノードが、検査すべきアプリケーションノードを、アプリケーションノードによって得られたリストを評価することによって求めるとよい。一構成例では本方法は、制御ノードにおいて実行される制御エージェントおよびアプリケーションノードにおいて実行されるアプリケーションエージェントによって実施される。

【0034】

他の有利な構成が従属請求項からもたらされる。更に、図面を参照しながら実施例に基づいて本発明を詳細に説明する。図1はクラスタ装置の第1の実施例を示し、図2は使用される概念を説明するためのダイアグラムを示し、図3はアプリケーションエージェントの機能概要を示し、図4は制御エージェントを示し、図5は図1による本発明装置の一部を示し、図6はアプリケーションエージェントおよび制御エージェントの一動作態様の実施例を示し、図7はプール形成を有するクラスタ装置の第2の実施例を示し、図8は読み書き可能なメモリ領域の一部を概略的に示し、図9はアプリケーションエージェントによって管理されるテストメントの一部を概略的に示す。

【0035】

図1は6つのデータ処理システムを有する本発明によるクラスタ装置を示す。サーバとして構成されているこれらのデータ処理システムのうちの4つは、アプリケーションノ

10

20

30

40

50

ド A K 1 , A K , A K および A K ' を成す。他の 2 つのサーバはそれぞれ 1 つの制御ノード K K を成す。全てのサーバは主プロセッサならびにメインメモリおよび/またはハードディスク記憶装置 S P 1 を有する。各アプリケーションノードもしくは制御ノードのメモリはオペレーティングシステム O S を有する。オペレーティングシステム O S は、アプリケーションノードにおけるプログラム制御のための機能、動作監視および維持のための機能、そしてノードの個々の構成要素へのアクセスのための機能を有する。更に、アプリケーションノード A K 1 , A K および A K ' のメモリ S P 1 は、それぞれ、オペレーティングシステム O S の一部であるアプリケーションエージェント A A を有する。制御ノード K K のメモリは制御エージェント K A を含む。

【 0 0 3 6 】

各アプリケーションノード A K 1 , A K および A K ' は第 1 のネットワーク N W 1 を介して制御ノード K K ならびに記憶装置 S P に接続されている。このネットワークは、アプリケーションノード A K 1 , A K および A K ' 間相互のデータ転送、アプリケーションノード A K 1 , A K および A K ' と制御ノード K K との間のデータ転送、そしてアプリケーションノード、制御ノードおよび記憶装置 S P 間のデータ転送を可能にする。第 2 のネットワーク N W 2 はアプリケーションノードおよび制御ノードをクライアントコンピュータ C L に接続する。クライアントコンピュータはアプリケーションノードに照会するために構成されていて、照会は処理のためにアプリケーションノードに送られる。

【 0 0 3 7 】

記憶装置 S P は全てのアプリケーションノード A K 1 , A K および A K ' のためのオペレーティングシステム O S を有する。各アプリケーションノードのアプリケーションエージェント A A はこのオペレーティングシステム O S のサービスであり、オペレーティングシステムの初期化後に始動される。それはバックグラウンドで動作するデーモンである。更に、記憶装置 S P は多数のプログラムモジュール M 1 , M 2 , M 3 を含む。これは、他方では個々のインスタンスに分割することのできる大きなアプリケーションである。例えばモジュール M 1 は 5 つのインスタンスを有し、モジュール M 2 は 2 つのインスタンスを有し、そしてモジュール M 3 は 1 つのインスタンス I からなる。

【 0 0 3 8 】

種々のモジュールのインスタンスがアプリケーションノード A K 1 および A K のメモリ S P 1 にロードされていて、そこで実行される。例えばアプリケーションノード A K 1 はモジュール M 2 のインスタンス I 1 ならびにモジュール M 1 のインスタンス I 2 および I 3 を実行し、両アプリケーションノード A K はインスタンス I 4 ~ I 7 ならびにインスタンス I を実行する。アプリケーション A K ' では他のインスタンスは実行されない。

【 0 0 3 9 】

モジュールとインスタンスとの間の関係は図 2 から読み取ることができる。モジュール M は複数のコンピュータで実行される比較的大きなアプリケーションである。このためにモジュール M はインスタンスと呼ばれる複数の小さな単位に分割されている。個々のインスタンスは、その必要が生じれば、互いに連絡し合っデータ交換する。それによって関連性が生じる。他方では、インスタンスは個々のプロセス P r からなるサブインスタンス S u I に分割されている。個々のインスタンスはそれぞれのサブインスタンス S u I およびプロセス P r と一緒に 1 つのコンピュータ上で実行される。この場合にノードは異なるモジュールまたは同じモジュールに属する複数のインスタンスを実行することもできる。例えば、インスタンス I 1 および I 4 を有するモジュール M 2 がアプリケーション A K 1 および A K に分配されている。ノード A K 1 上では同時にモジュール M 1 のインスタンス I 2 および I 3 が実行される。

【 0 0 4 0 】

各インスタンスには一義的な識別番号 I P 1 , I P 2 および I P 3 が割り当てられていて、アプリケーションノードにおいてインスタンスを実行する際には、そのアプリケーションノードにこの識別番号が割り振られる。インスタンスが実行されるアプリケーションノードは、この一義的な識別番号 I P 1 , I P 2 および I P 3 を確認することができる。

10

20

30

40

50

したがって、あるアプリケーションノードから他のアプリケーションノードへのインスタンスの交替が問題なく可能である。なぜならば、識別番号が古いアプリケーションノードでは抹消され、新たなアプリケーションノードに割り振られるからである。インスタンスへのアクセスは該当する一義的な識別番号を用いて行なわれる。本実施例では一義的な識別番号は仮想的なIPアドレスによって定義されている。したがって、アプリケーションノードAK1は、インスタンスの仮想的なIPアドレスIP1, IP2およびIP3を受け取る。ネットワークNW2を介してインスタンスI1にアクセスしようとするクライアントCLは、インスタンスI1に割り付けられている仮想的なIPアドレスに照会を送信する。アプリケーションノードAK1はこの照会を受信してインスタンスI1に転送し、それをインスタンスI1が処理する。インスタンスに割り付けられていてそのインスタンスの実行時にアプリケーションノードに割り当てられる仮想的なアドレスの使用は、アプリケーションノードの自由な選択を可能にする。

10

【0041】

アプリケーションエージェントAAの種々の機能が図3においてアプリケーションノードAK1の例で説明されている。このアプリケーションノードではインスタンスI1, I2およびI3が実行される。アプリケーションエージェントAAは監視手段UBを有し、それによりインスタンスを監視する。これに属するのが、例えば各インスタンスのための使用メモリのプロセッサ稼働率の測定、処理された照会および他の動作パラメータである。更に、アプリケーションエージェントAAは正しい動作態様およびインスタンスI1およびI3の可用性を監視する。更に、アプリケーションエージェントは、監視手段により、アプリケーションノードにおいて他の監視すべきインスタンスが存在するかどうかを検査する。監視手段は、監視すべきであるインスタンスを認識するように構成されている。これに加えて、監視手段は、アプリケーションノードにおいて実行されるインスタンス間の関連性を認識する。さらに、監視手段はとりわけノード上で実行中の全てのプロセスのリストを周期的に分析する。周期的な検査によって後から始動されたインスタンスも認識されて自動監視に引き継がれる。

20

【0042】

更に、アプリケーションエージェントは、以後の経過においてテストメントと呼ばれるリストTの発生Lのための機能を有する。このテストメントTは、個別の監視すべきインスタンスI1~I3の全ての重要なデータDが割り付けられた個々の部分テストメントTI1~TI3からなる。割り付けられたデータDには、インスタンスの名称のほかにはインスタンスの誤りのない動作のために必要な動作パラメータも属している。これのための例が、必要なメモリおよび計算容量、環境変数、他のインスタンスおよび動作パラメータに対するインスタンス同士の関連性などである。付加的にテストメントにはアプリケーションノードに関するデータおよびパラメータが含まれている。これらは、例えば使用されるサーバの型および種類、名称、位置、メモリおよびプロセッサである。全てのアプリケーションノードのテストメントのこれらのパラメータの評価がクラスタ構造の決定を可能にし、かつ他の設定可能性を与える。

30

【0043】

監視機能が監視すべきインスタンスを見つけ出さない場合に、エージェントは、アプリケーションノードが新たなインスタンスの実行のために使用されることを認識し、このことを相応にテストメントにおいてはっきり示す。

40

【0044】

アプリケーションエージェントAAは始動機能STにより構成されている。したがって、アプリケーションエージェントAAは、記憶装置SPからネットワークNW1を介してインスタンスを取り込み、これをアプリケーションノード上で実行する。実行されるインスタンスは監視手段UBによって周期的に検査される。

【0045】

評価および判定機能BEによりアプリケーションエージェントAAは新たなインスタンスが実行可能であるかどうかを評価し、アプリケーションエージェントAAはインスタ

50

スが始動されるべきであるアプリケーションノードを的確にとらえる。アプリケーションノードの動作パラメータ（プロセッサおよびメモリ稼働率）の測定および内部の記憶（テストメントTはこの記憶の一部である。）のほかに、新たなインスタンスの始動に関する判定は定められた規則に依存する。条件が満たされたならばアプリケーションエージェントAAはインスタンスを取り込んでこれを実行する。

【0046】

定められた規則の例は、例えばプロセッサ能力およびメモリについての最小限準備の条件である。他の規則は、特定のインスタンスを定められた時間の間のみ実行するという定義付けである。ここでも他の規則は、新たなインスタンスの始動時に監視機能によりインスタンスの関連性を調べ、このインスタンスに関連したこれまで始動されていないインスタンスを同様に実行に至らしめることを意味する。

10

【0047】

新たなインスタンスの始動後にアプリケーションエージェントはメッセージMEをネットワークNW1を介して別のアプリケーションノードの他のアプリケーションエージェントに送信する。これにより、新たなインスタンスの始動が成功したことを示す。

【0048】

アプリケーションエージェントは、個々のインスタンスI1, I2およびI3のための監視手段UBによって、監視されるインスタンスI3の不慮の予定より早い終了を認識することができる。更に、誤りのない動作を維持するために、アプリケーションエージェントは障害のあるインスタンスI3の終了および再始動のための機能を有する。再始動が成功しなかった場合に、エージェントは、テストメントTから、障害のあるインスタンスI3の部分テストメントTI3を発生させ、部分テストメントTI3を有するこのインスタンスの始動のための要求AFをネットワークNW1を介して他のアプリケーションノードに送信する。その際にエージェントは障害のあるインスタンスI3の終了または更なる始動試行の停止を行なうことができる。

20

【0049】

更に、アプリケーションエージェントはアプリケーションノードにおいて実行されるインスタンスを終了させる機能を有する。この機能はインスタンス終了要求にしたがって使用される。それによってインスタンスが1つのアプリケーションノードにおいて終了させられ、他のアプリケーションノードにおいて新たに実行される。

30

【0050】

アプリケーションノードAKにおけるアプリケーションエージェントAAの独立かつ自立した判定到達によって、制御ノードまたは中央で動作する監視手段によるアプリケーションノードAKにおける個々のインスタンスの連続的な制御および監視はもはや必要でない。

【0051】

図4は制御エージェントKAの機能に関する概要を示す。制御エージェントKAはアプリケーションエージェントAAからそれぞれのテストメントを受け取り、これらを管理する。これによりアプリケーションエージェントがクラスタに登録される。制御エージェントは、アプリケーションエージェントAAのテストメントTから、クラスタ内に存在する全てのアプリケーションノードAKをそれらのハードウェア情報を含めて備えたリストを発生する。それによって、制御エージェントは自立的にクラスタの現在の構成情報を入手しかつ動的な変化も登録する。更に、制御エージェントKAは、ネットワークNW1を介する全てのアプリケーションノードAKの機能性および存在の検査のための手段UPを備えている。アプリケーションノードAKの機能性および存在は、アプリケーションエージェントからの簡単な存在信号の送信によって伝達される。例えば、制御エージェントKAがネットワークNW1を介してPing信号を個々のアプリケーションノードAKに送信することができる。

40

【0052】

機能テストにおける応答不在によって知らされるアプリケーションノード障害の際には

50

、制御エージェント K A は当該アプリケーションノード A K に対するテストメントを評価し、それから部分テストメント T I 3 を抽出する。この部分テストメントはこのインスタンスの実行のための要求 A F と共にネットワーク N W 1 に導かれ、そして残っているアプリケーションノード A K に導かれる。これの代替として、制御エージェント K A は、インスタンスを実行するアプリケーションノードを決定するための機能を持っている。図 1 のクラスタ装置においてアプリケーションノード A K ' はインスタンスを持っていないので、制御エージェント K A は、アプリケーションノード A K 1 の障害後におけるインスタンス I 1 , I 2 および I 3 の実行のために、このアプリケーションノード A K ' を決定する。アプリケーションノード A K 1 のアプリケーションエージェント A A から伝達されるテストメント T によって、中止されたインスタンスがアプリケーションノード A K ' において始動可能となる。 10

【 0 0 5 3 】

図 5 は新しいアプリケーションノード A K " が付け加えられた本発明によるクラスタ装置の部分図を示す。アプリケーションノード A K では 2 つのインスタンス I 1 および I 2 が実行される。アプリケーションノード A K ' ではインスタンスは全く実行されない。アプリケーションノード A K のアプリケーションエージェント A A は、两部分テストメント T I 1 および T I 2 を有するテストメント T を作成し、これを制御ノード K K および制御エージェント K A に伝達したところである。ノード A K ' のアプリケーションエージェントは空のテストメント T ' を制御エージェント K A に伝達し、そのテストメント T ' における登録 S P によりアプリケーションノード A K ' がインスタンスの始動のために準備完了であることを知らせる。この登録によりノード A K ' が自由なノードとして明らかにされる。 20

【 0 0 5 4 】

制御エージェント K A は、自身の側で、アプリケーションノード A K および A K ' のテストメント T および T ' を有するリストを管理する。エージェント K A は、周期的にノードのアプリケーションエージェント A A の状態信号を要求することによって、ノードの存在を検査する。アプリケーションノード A K における監視されるインスタンス I の動作パラメータの変化、すなわち終了または新たなインスタンスの始動の際には、この変化がそれぞれのアプリケーションエージェント A A によって自動的に制御ノード K K の制御エージェント K A に伝達される。したがって、制御エージェント K A のリストは、常にアプリケーションノード A K のテストメントの現在状態を含んでいる。更に、制御エージェントはアプリケーションノードのハードウェアパラメータに関する情報を受け取る。 30

【 0 0 5 5 】

ここで、新しいアプリケーションノード A K " がネットワーク N W 1 に接続される。初期化段階後にアプリケーションエージェント A A がノード A K " において始動する。エージェント A A の監視機能 U B が、アプリケーション A K " において実行されるプロセス、インスタンスおよびアプリケーションを検査し、自動的にアプリケーションエージェント A A によって監視すべきインスタンス I 3 を認識する。アプリケーションノードの動作パラメータと一緒に、エージェントはそれからインスタンス I 3 の動作に必要な全てのデータおよびパラメータを含んだ部分テストメントを有するテストメント T " を発生する。アプリケーションノード A K " の発生させられたテストメント T " は制御エージェント K A に伝達される。それによりアプリケーションエージェントがクラスタにおいて登録され、アプリケーションノード A K " においてインスタンスが実行され、そのインスタンスが監視される。制御エージェント K A は今やアプリケーションノード A K , A K ' および A K " の存在をそのリストにあるテストメントにしたがって検査する。 40

【 0 0 5 6 】

アプリケーションノード A K " がネットワーク N W 1 から分離されるか、または予定より早く例えば電源障害によって遮断された場合には、存在の検査が否定的結果をもたらす。制御エージェント K A はテストメントにあるインスタンスの実行要求を有するテストメント T " をアプリケーションノード A K および A K ' に送信する。アプリケーションエー 50

ジェント A A はそのテストメントを受信し、測定、記憶および外部パラメータにより、全体のテストメントまたは部分テストメントがアプリケーションノードにおいて実行可能であるか否かの判定が的確にとらえられる。

【 0 0 5 7 】

アプリケーションノード A K ' のアプリケーションエージェントは肯定的判定を的確にとらえて全体のテストメント T " を受け取る。エージェントはインスタンス I 3 をテストメントにおいて予め与えられたパラメータにしたがってそのノードにて新たに始動し、今や新たなインスタンス I 3 の部分テストメントを含む新しいテストメント T ' を制御エージェントに伝達する。

【 0 0 5 8 】

図 6 は他の好ましい構成を示す。クラスタ装置は 2 つのアプリケーションノード A K および A K ' と記憶装置 S P と制御ノード K K とを有し、これらはネットワーク N W 1 を介して互いに接続されている。ノード相互および記憶装置 S P との通信は T C P / I P プロトコルにより行なわれる。

【 0 0 5 9 】

実行されかつアプリケーションエージェントによって監視されるインスタンス I 1 , I 3 および I 2 は記憶装置 S P に保存されているモジュール M 2 を成す。更に記憶装置 S P は、モジュール M 1 と、アプリケーションノードに共通に使用されるオペレーティングシステム O S とを含み、オペレーティングシステムはアプリケーションエージェント A A を有する。記憶装置 S P は、2 つの部分領域 B 1 および B 2 に分割されている領域 B を有する。

【 0 0 6 0 】

領域 B は、全てのアプリケーションエージェント A A および制御エージェント K A のために読み書き可能に構成されている。アプリケーションエージェントは、それらのアプリケーションノードのテストメントを記憶装置 S P における部分領域 B 1 に保存する。アプリケーションノードにおける変化の際には、このノードのアプリケーションエージェントが新たなテストメントを発生し、それにより領域 B 1 において古いテストメントを交換する。制御ノードの制御エージェントは領域 B 1 におけるテストメントを評価し、それにもなって監視すべきアプリケーションノードのリストを発生する。

【 0 0 6 1 】

更に、各ノードのアプリケーションエージェントは、記憶装置の領域 B 2 を周期的に評価する。領域 B 2 にはインスタンスの実行要求が保存されている。この構成においては、要求は領域 B 2 へのテストメントまたは部分テストメントの格納によって行なわれる。アプリケーションエージェントは、領域 B 2 に格納されているテストメントまたは部分テストメントを読み取って、実行に関する独立した判定を的確にとらえる。アプリケーションノードがテストメントを引き受けることができる場合には、アプリケーションエージェントはそのテストメントを領域 B 2 から消去し、指定されたインスタンスを始動する。要求または引き受け後の通知は、領域 B 2 へのテストメントの格納またはその領域からのテストメントの消去によって簡単なやり方にて行なわれる。要求の拒否はテストメントが領域 B 2 に残されていることによって自動的にもたらされる。

【 0 0 6 2 】

インスタンスを終了させるアプリケーションエージェントはそのテストメントを領域 B 2 に格納するので、他のアプリケーションノードがこれを引き受けることができる。アプリケーションノードが完全に故障し、アプリケーションエージェントがこれを前もって領域 B 2 へのテストメントの格納によって指定することができない場合には、制御エージェントが故障したアプリケーションノードのテストメントを領域 B 2 へ移動する。残りのノードのアプリケーションエージェントはそれらの側で判定を的確にとらえる。このやり方で高い柔軟性が達成される。多数の部分テストメントへのテストメントの分割によって、障害のあるアプリケーションノードのインスタンスを複数のノードに配分することができる。アプリケーションエージェントの独立および装置 S P の共通使用されるメモリによっ

10

20

30

40

50

て、制御ノード K K の障害発生時にも誤りのない動作が保証されている。

【 0 0 6 3 】

ここに挙げた例のほかに、多数の他の構成を見いだすことができる。特にアプリケーションエージェントの判定到達のための規則、制御エージェントおよびアプリケーションエージェントの機能および課題、そしてテストメントにおけるパラメータが拡張可能である。

【 0 0 6 4 】

本発明の他の観点は、いわゆるクラスタ装置内の個別データ処理システムの仮想クラスタへのグループ化に関する。この場合に、クラスタ装置内の幾つかのデータ処理システムに同一識別番号が割り当てられ、そのようにしてこれらのデータ処理システムが1つのプールの統合される。ここにおいて、仮想という概念は、異なるデータ処理システム相互における規則によって定められた論理的な関連性にすぎない。複数のデータ処理システムによるクラスタ装置内のプール形成は、異なるアプリケーションを高可用性に保とうとする場合に特に有利である。幾つかのデータ処理システムを特別にデータベースサービスの実行のために設け、これに対して同じ物理的なクラスタの他のデータ処理システムはウェブアプリケーションのために設けることが望ましい。

【 0 0 6 5 】

プール形成は一般的な規則にしたがって行なわれる。これらの規則は、例えば高可用性のアプリケーションに対する規則に関連し得るが、しかし例えば純粋なハードウェアパラメータも含み得る。更に、物理的なクラスタ内でのプール形成により、異なるユーザグループの個別プールを割り振ることが可能である。それぞれのユーザグループによって始動されたアプリケーションは、それぞれのプールに割り当てられたデータ処理システムにおいてのみ実行されて、高可用性に保たれる。相互に割り当てられている複数のデータ処理システムからなるこのようなプールは、物理的なクラスタ装置内の仮想クラスタとも呼ばれる。

【 0 0 6 6 】

図 7 は、その中に含まれている複数の仮想クラスタもしくはプールを有するクラスタ装置の実施形態を示す。そこに示されたクラスタは、共通のネットワーク N W 1 を介して複数のデータ処理システムに接続されている複数の記憶装置 S P , S P 1 および S P 2 を含んでいる。これらのデータ処理システムはそれぞれコンピュータとして主プロセッサおよび主メモリを装備している。これらのデータ処理システムのうち 1 2 個がアプリケーションノード A K 1 ~ A K 1 2 として構成されている。他の 2 つのデータ処理システムが制御ノード K K および K K 1 を成す。アプリケーションノード A K 1 ~ A K 1 2 ではそれぞれ 1 つのアプリケーションエージェント A A が実行される。制御ノード K K および K K 1 はそれぞれ 1 つの制御エージェント K A を含んでいる。

【 0 0 6 7 】

物理的なクラスタ装置は、この実施例において、3つの仮想クラスタ V C 1 , V C 2 および V C 3 を含む。仮想クラスタ V C 1 は、3つのアプリケーションノード A K 1 0 , A K 1 1 および A K 1 2 と、記憶装置 S P 1 とを含んでいる。仮想クラスタ V C 2 は、アプリケーションノード A K 1 ~ K 5 と、制御ノード K K 1 と、記憶装置 S P 2 とを含んでいる。更に、アプリケーションノード A K 3 および A K 4 は、アプリケーションノード A K 6 ~ A K 9 と一緒に仮想クラスタ V C 3 に割り当てられている。制御ノード K K および記憶装置 S P はこの実施例では仮想クラスタの一部ではない。

【 0 0 6 8 】

仮想クラスタのそれぞれのプールへの個々のアプリケーションノード A K もしくは制御ノード K K の割り振りは一般的な規則を介して行なわれる。これらの規則は、部分的には、外部のユーザによって予め与えられるか、制御ノード K K の制御エージェント K A によって予め与えられるか、またはアプリケーションノード A K における個々のアプリケーションエージェントのテストメントからもたらされる。個々の仮想クラスタの同定のために、そして 1 つのプールへの物理的なクラスタの個々のノードの割り振りのために、I P アド

10

20

30

40

50

レスの一部またはIPアドレス自体を使用するのが有利である。例えば、仮想クラスタVC1のアプリケーションノードAK10, AK11およびAK12には、部分的に一致するIPアドレスが割り付けられる。IPアドレスが同じ部分を有する物理的クラスタ装置のノードは、同じプールもしくは仮想クラスタに属する。この仮想クラスタにおけるノードからまたはノードへの通信は同様にこの識別番号を含んでいる。データ処理システムにおけるエージェントの相応の評価によって他の識別番号を有する通信は無視される。

【0069】

仮想クラスタVC2のアプリケーションノードAK1においては、アプリケーションエージェントAAならびにインスタンスL1およびインスタンスI1が実行される。アプリケーションノードAK2はインスタンスL2およびI2を含む。アプリケーションノードAK1のインスタンスI1およびアプリケーションノードAK2のインスタンスI2は、共通に1つのアプリケーションを成す。これらの共通なアプリケーションは、仮想クラスタVC2の異なるアプリケーションノードAK1~AK5において高可用性に保たれている。したがって、共通アプリケーションIのための両インスタンスの一方I1もしくはI2の始動のための要求は、相応のアプリケーションノードが仮想クラスタVC2に割り当てられているときにのみ引き受けられる。

【0070】

アプリケーションノードAK3は、アプリケーションエージェントAAのほか部分インスタンスL3ならびにインスタンスL31およびL32も含み、これらはそれぞれ高可用性にて実行される。インスタンスL3はアプリケーションノードAK2およびAK1のインスタンスL2およびL1と一緒に仮想クラスタVC2の他のアプリケーションを成す。アプリケーションノードAK4およびAK5は予備ノードであり、予備ノードにおいては仮想ノードVC2の他のインスタンスは実行されない。

【0071】

更に、アプリケーションノードAK3およびAK4は、仮想クラスタVC3の構成要素でもある。したがって、要求の評価およびインスタンス始動要求の発送のために、アプリケーションノードAK3およびAK4におけるアプリケーションエージェントAAが相応の要求を常に同じプールに所属しているノードのアプリケーションエージェントに送信することが必要である。このために、例えばノードAK3におけるアプリケーションエージェントAAが、仮想クラスタへの個々のインスタンスの一義的な割り当ての相関性に関して拡張されたテストメントを既に含んでいる。

【0072】

このテストメントからの一部が図9において見ることができる。テストメントは2つの比較的大きな部分領域に分けられていて、部分領域はそれぞれノードAK3を構成部分とする両プールのアプリケーションを含む。これはプールもしくは仮想クラスタVC2ならびにVC3である。仮想クラスタVC2は、仮想クラスタへの一義的割り当てを可能にする識別番号IP1を含む。更に、インスタンスL3が仮想クラスタVC3において実行される。したがって、インスタンスL3に関係する全ての通知に対して対応関係IP1が一緒に送られる。ノードが同じ対応関係を持っていないところのエージェントは、この仮想クラスタの一部ではなく、したがって通知を無視する。

【0073】

第2の部分領域は、仮想クラスタVC3に割り当てられかつノードで実行されるインスタンスに対する全ての部分テストメントを含む。プールVC3は識別番号IP3を有する。したがって、テストメントのこの部分領域内においては、アプリケーションエージェントAAによってインスタンスL31およびL32が管理される。それぞれの部分テストメントL31およびL32は、これらの部分テストメントを実行するためのパラメータのほか、仮想クラスタVC3への対応関係も含んでいる。インスタンスの1つの故障および不成功に終わった故障インスタンスの新たな初期化の際に、ノードAK3のアプリケーションエージェントAAは、このインスタンスが実行される仮想クラスタに対する割り当て識別番号を有する部分インスタンスを発生する。インスタンスL31およびL32の障害

10

20

30

40

50

時には、仮想クラスタV C 3に識別番号I P 3によって割り当てられている部分インスタンスが発生され、インスタンスL 3の障害時には割り当て識別番号I P 2を有する相応の部分インスタンスが発生される。

【0074】

エージェントA Aは、アプリケーションノードに、このインスタンスの実行要求を送信する。アプリケーションノードにおけるそれぞれのアプリケーションエージェントはこの要求を評価して、手始めにそれが同じ識別を有する仮想クラスタの部分であるか否かを検査する。それがインスタンスを実行させるべき仮想クラスタの部分でない場合には、実行要求は無視される。ほかの場合にはそれに必要な実行のためのリソースが使用可能であるか否かが検査される。

10

【0075】

プールV C 2には付加的に、可用性および障害安全性を高めるために、制御エージェントK Aを有する制御ノードK K 1が割り当てられている。このために、アプリケーションノードA K 1 ~ A K 5におけるアプリケーションエージェントA Aが、それらの部分インスタンスを発生し、それらを大容量記憶装置S P 2における共通な読み書き可能なメモリ領域に格納する。制御ノードK K 1における制御エージェントK Aは、個々のアプリケーションノードA K 1およびそれらのエージェントの機能性を規則的な時間間隔で送出される状態メッセージによって監視する。

【0076】

このメモリ領域およびアプリケーションノードA K 1 ~ A K 5におけるアプリケーションエージェントの個々の伝達されるテストメントの概略図を図8に示す。全体リストにおいては、とりわけどのアプリケーションノードA Kが仮想クラスタV C 2に割り当てられているかが整理されている。更に、仮想クラスタV C 2において目下のところ実行されるアプリケーションを有するリストが作成される。詳細には、これは部分インスタンスI 1およびI 2を有するアプリケーションならびに部分インスタンスL 1, L 2およびL 3を有するアプリケーションである。これらのアプリケーションのそれぞれに優先順位が割り付けられている。それぞれの部分インスタンスはこれらの優先順位を受け継ぐ。優先順位は、仮想クラスタV C 2において実行される個々のアプリケーションがどの程度重要であるかを指定する。したがって、優先順位は実行されるアプリケーションの順序もしくはランク順を成す。

20

30

【0077】

この例では、两部分インスタンスI 1およびI 2を有するアプリケーションが優先順位指標P S 1を持ち、部分インスタンスL 1, L 2およびL 3を有するアプリケーションが優先順位指標P S 2を持つ。この場合に優先順位指標P S 2は指標P S 1よりも小さい。したがって、部分インスタンスL 1, L 2およびL 3を有するアプリケーションが部分インスタンスI 1およびI 2を有するアプリケーションよりも重要でない。

【0078】

更に、記憶装置S P 2のメモリ領域は個々のアプリケーションノードA K 1 ~ A K 5のテストメントT 1 ~ T 5を含んでいる。これらは、それぞれのアプリケーションノードA K 1 ~ A K 5において作動する部分インスタンスのための部分テストメントを含んでいる。アプリケーションノードA K 4およびA K 5のテストメントT 4およびT 5は空である。

40

【0079】

制御ノードK Kの制御エージェントK Aは、一般に個々のアプリケーションノードの高可用性を監視する。今、例えばアプリケーションノードA K 1が完全に故障した場合には、もはや部分インスタンスI 1もしくはL 1も実行されない。制御ノードK Kにおける制御エージェントK Aは、今やテストメントT 1から部分インスタンスI 1およびL 1のための2つの部分テストメントを発生する。その際に高いほうの優先順位指標に基づいてインスタンスI 1を有する部分テストメントがこの部分インスタンスの実行要求と一緒に物理的なクラスタ内における個々のアプリケーションノードにネットワークを介して送信さ

50

れる。この部分テストメント内において、実行すべきインスタンス I 1 がどの仮想クラスタに割り当てられているかが指定されている。

【0080】

仮想クラスタ VC 2 に割り当てられていないアプリケーションノード AK におけるアプリケーションエージェント AA は、実行要求を無視する。それに反してアプリケーションノード AK 2 ~ AK 5 におけるエージェント AA はそれらのリソースを検査する。アプリケーションエージェント AA の 1 つが、場合によっては部分テストメントを引き受け、インスタンス I 1 を自身のノードにおいて実行へと至らしめる。部分インスタンス I 2 の始動後に相応のメッセージが制御エージェント KA に返送される。部分インスタンス I 1 が部分テストメントと一緒に仮想クラスタ内のアプリケーションエージェント AA の 1 つによって引き継がれて実行に成功したときにはじめて、制御エージェント KA が部分インスタンス L 2 を有する部分テストメントを実行要求と一緒に送信する。

10

【0081】

個々のアプリケーションもしくはインスタンスの優先順位付与によって、高い優先順位を有するアプリケーションが常に高可用性に保たれる。十分な容量が存在するときのみ、低い優先順位を有するインスタンスも新たに実行に至らしめられる。この実施例においては、仮想クラスタ VC 2 内の制御ノード KK 1 における制御エージェント KA が、部分テストメントの発生および部分インスタンスの実行要求を引き受ける。

【0082】

アプリケーションノード AK 3 の障害の見極めがつく場合には、これをノード AK 3 のアプリケーションエージェント AA が確認する。このノードのエージェント AA は部分インスタンス L 3 を有する部分テストメントおよび実行要求を発生し、これを物理的クラスタ装置および仮想クラスタ装置のアプリケーションノードにおけるエージェントに送信する。更に、アプリケーションノード AK 3 におけるアプリケーションエージェント AA は、部分インスタンス L 3 1 および L 3 2 を有する 2 つの部分テストメントを発生し、これらを同様に実行要求と共にクラスタ装置に送信する。しかしながら、部分インスタンス L 3 1 および L 3 2 は仮想クラスタ VC 3 に割り当てられていて、アプリケーションノード AK 1 , AK 2 および AK 5 によって無視される。しかしながら、適切な自由なリソースにおいては、アプリケーションノード AK 4 もしくは AK 6 ~ AK 9 が部分インスタンス L 3 1 および L 3 2 を引き受けることができる。

20

30

【0083】

この実施例においては、メッセージが物理的なクラスタ装置内の全てのアプリケーションノードに送信される。しかしながら、メッセージが同じプールにおけるノードに由来する場合にのみ処理が行なわれる。拡張においては、同じ仮想クラスタ内のノードのみにメッセージを送信することもできる。それにより、確かにデータ発生が低減されるが、しかし柔軟性も縮小される。

【0084】

加えて、アプリケーションノードにおいて計画的でなく終了させられる低い優先順位のインスタンスは、次の場合にはアプリケーションノード全体の再始動のためにアプリケーションエージェントを動かすことができないように配慮することが望ましい。すなわち、このノードにおいて、より高い優先順位を有するインスタンスがなおも誤りなく実行される場合である。例えば、アプリケーションノード AK 2 におけるアプリケーションエージェント AA は、部分インスタンス L 2 の障害発生時に、より高い優先順位を有する部分インスタンス I 2 がなおも誤りなく実行される場合には、アプリケーションノード AK 2 の全体の完全な再始動を行なわない。したがって、再始動のためには、アプリケーションエージェント AA が、部分インスタンス I 2 を有する部分テストメントおよび実行要求を、仮想クラスタ VC 2 のアプリケーションノードに送信しなければならない。この部分インスタンスの引き受けおよび実行成功の確認時に、アプリケーションノード AK 2 におけるアプリケーションエージェント AA は、アプリケーションノード AK 2 の完全な再始動を初期化する。

40

50

【 0 0 8 5 】

部分IPアドレスの助けによる個々の仮想クラスタへの割り当てによって、非常に動的にかつ柔軟にリソース要求における有り得る変化に対応することができる。付加的に仮想クラスタ内においても個々のデータ処理システム間における更なるグループ化またはプール形成を設定することができる。仮想クラスタVC3においては、例えばアプリケーションノードAK8およびAK9が仮想クラスタ内において他のグループを成す。このグループ化も一般的な規則を介して制御可能である。更に、完全なクラスタ構造を監視しかつ個々のデータ処理システムを規則的間隔にて存在を監視する制御エージェントKAを有する他の制御ノードKKが設けられるとよい。他のデータ処理システムが物理的クラスタに付け加えられるならば、この制御ノードは、仮想容量増大のために、付け加えられたデータ処理システムを異なる仮想クラスタに割り当てることができる。識別番号の割り付けによるプール形成の導入ならびにプール内の個々の部分インスタンスの優先順位付与は、物理的クラスタ装置内における非常に細かい段階付けおよび選択制御を可能にする。この場合に、個々のアプリケーションノードならびに物理的クラスタ装置内におけるプールが大幅に自動的に構成可能である。付加的に管理上の課題が大々的に解消する。

10

【 0 0 8 6 】

以上のとおり、制御ノードとアプリケーションノードとの間の明白に少ないデータ交換がこれと同時に大きな柔軟性をともなって達成される。個々のアプリケーションノードの自立によって完全に制御ノードを省略することさえも可能である。ノード経過の管理、制御および監視はアプリケーションエージェントの務めであり、これらのアプリケーションエージェントは互いに直接の通信を持っていない。アプリケーションエージェントの形成はインスタンスの独立の認識および監視を可能にする。それによって、高コストの構成を省略することができ、クラスタ構造が自立的に発生させられるためにクラスタ構造に関する正確な知識は必要でない。しばしばノード数を変更するクラスタを使用する場合にはまさに、独立した監視のこの構想は高い柔軟性をもたらす。

20

【 図面の簡単な説明 】

【 0 0 8 7 】

- 【 図 1 】 クラスタ装置の第 1 の実施例を示す概略図
- 【 図 2 】 使用される概念を説明するための概略図
- 【 図 3 】 アプリケーションエージェントの機能の概略図
- 【 図 4 】 制御エージェントを示す概略図
- 【 図 5 】 図 1 による本発明装置の部分詳細図
- 【 図 6 】 アプリケーションエージェントおよび制御エージェントの一動作態様の実施例を示す概略図
- 【 図 7 】 プール形成を有するクラスタ装置の第 2 の実施例を示す概略図
- 【 図 8 】 読み書き可能なメモリ領域の一部を示す概略図
- 【 図 9 】 アプリケーションエージェントによって管理されるテストメントの一部を示す概略図。

30

【 符号の説明 】

【 0 0 8 8 】

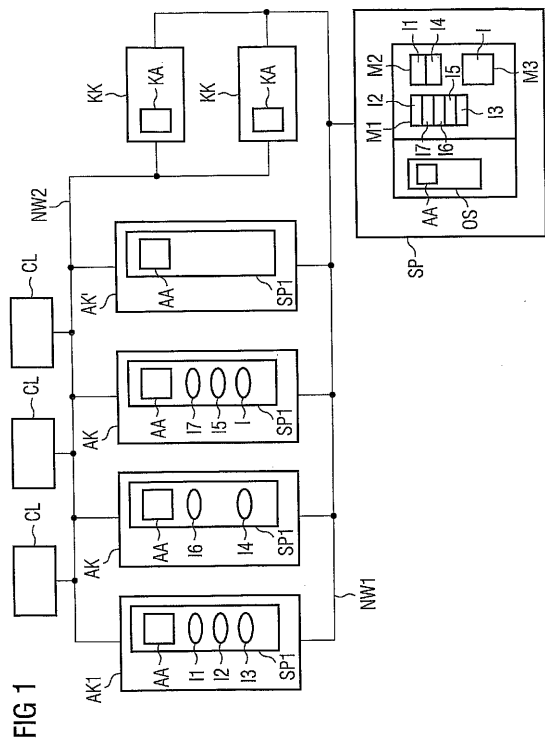
AK , AK , AK ' , AK "	アプリケーションノード
KK	制御ノード
KA	制御エージェント
AA	アプリケーションエージェント
CL	クライアントコンピュータ
NW 1 , NW 2	ネットワーク
SP , SP 1	記憶装置
M 1 , M 2 , M 3	モジュール
I 1 , I 2 , . . . , I 7	インスタンス
SUI	副インスタンス

40

50

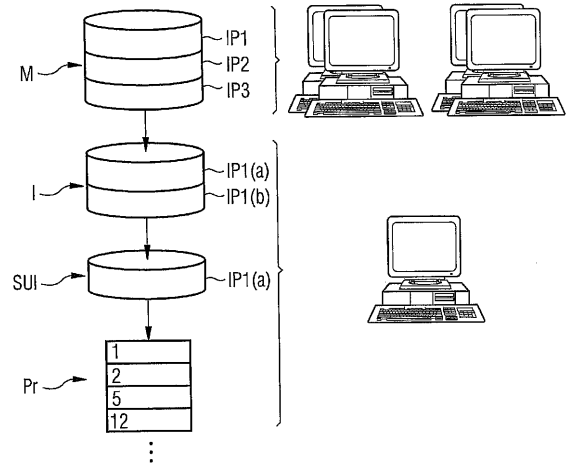
P r	プロセス
I P 1 , I P 2 , I P 3	識別番号
O S	オペレーティングシステム
T , T ' , T "	テストメント
T 1 , T 2 , T 3	部分テストメント
D	データ
L , U B , S T , M E , A F , B E	機能
U P , B S	機能
S P	識別番号
B , B 1 , B 2	メモリ領域

【 図 1 】



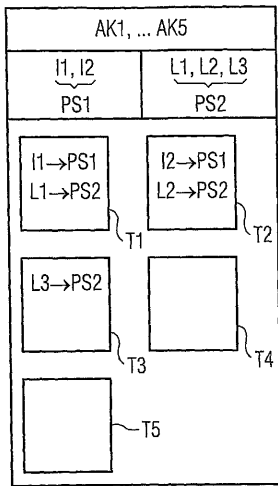
【 図 2 】

FIG 2



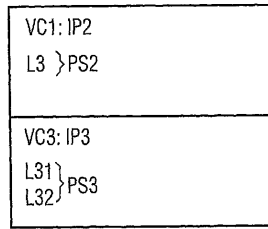
【 図 8 】

FIG 8



【 図 9 】

FIG 9



フロントページの続き

(81) 指定国 AP(BW, GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW

(72) 発明者 クラウス、ルディ

ドイツ連邦共和国 67158 エラーシュタット ポルトゥギーザーリング 53

Fターム(参考) 5B042 GA12 GA19 GC10 GC18 JJ03 JJ05 KK17

【要約の続き】

リケーションノードのアプリケーションエージェント(AA)へのインスタンス実行要求(AF)、

- ネットワーク(NW1)に接続されているアプリケーションノードのアプリケーションエージェント(AA)へのインスタンス実行要求(AF)引き受け後のメッセージ通知(ME)。